

Wouter Kool, Herke van Hoof, Max Welling

GUMBEL

Mathemagic



UNIVERSITY OF AMSTERDAM

ORTEC

OPTIMIZE YOUR WORLD

ANLAB

Amsterdam
Machine Learning Lab

- "This is how you
randomize a beam search!"

Wouter Kool, Herke van Hoof, Max Welling

- "You will *never* have
duplicate samples again!"

STOCHASTIC BEAMS

AND WHERE
TO FIND THEM

The Gumbel-Top- k Trick for Sampling
Sequences Without Replacement



 UNIVERSITY OF AMSTERDAM

ORTEC **ANLAB**
OPTIMIZE YOUR WORLD Amsterdam Machine Learning Lab



TL;DR

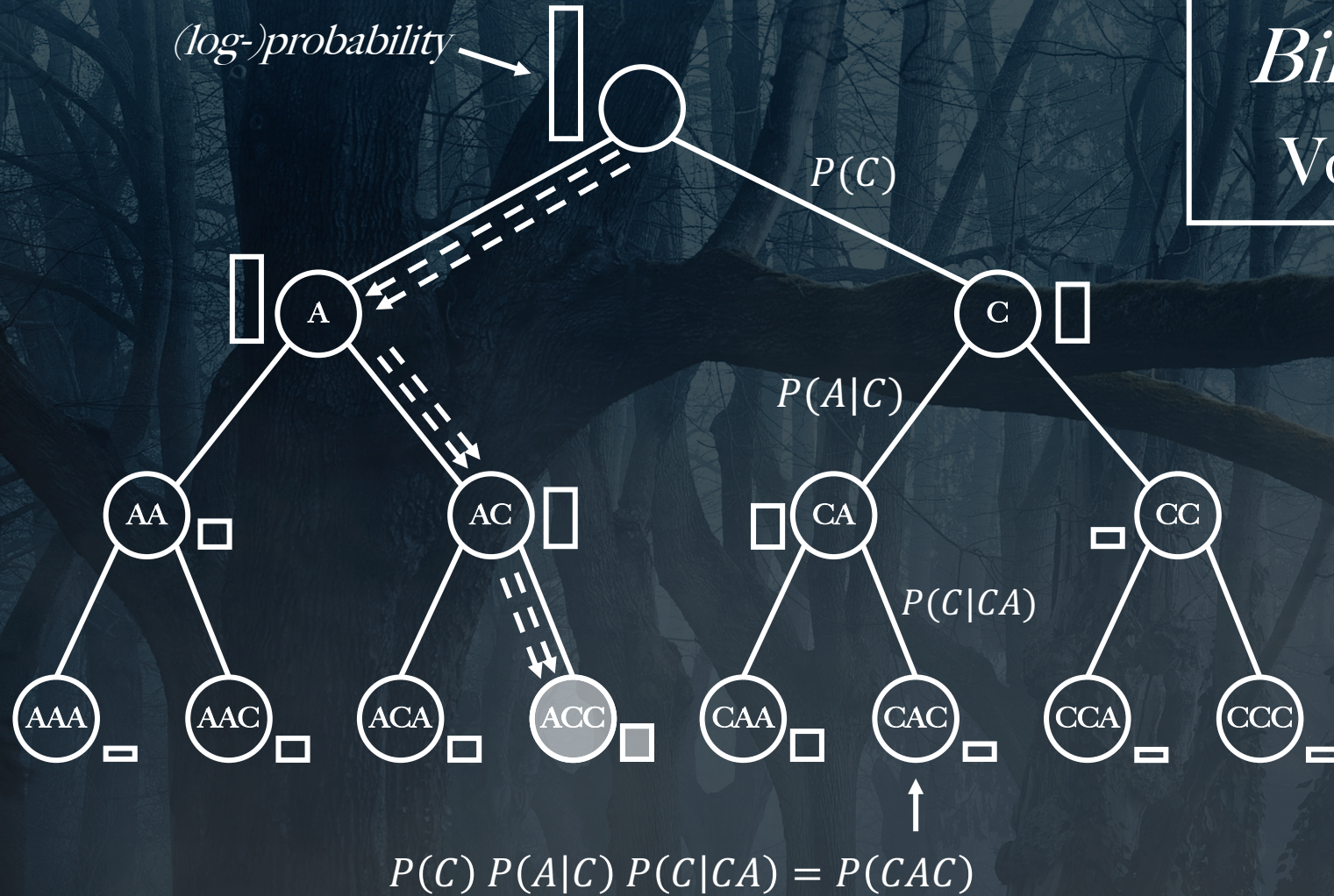
Stochastic Beam

Search finds a set of
unique samples
(without replacement)
from a sequence model.

Example

Binarese language model

Vocabulary: $\{A_{bra}, C_{adabra}\}$



*What if we want
a sample from
our model?*

The Gumbel-Max Trick

"Prof. Gumbeldore"

(Gumbel, 1945;
Maddison et al., 2014)



$$\phi_i = \log p_i$$

log-probability

$$G_i \sim \text{Gumbel}(0)$$

Gumbel noise

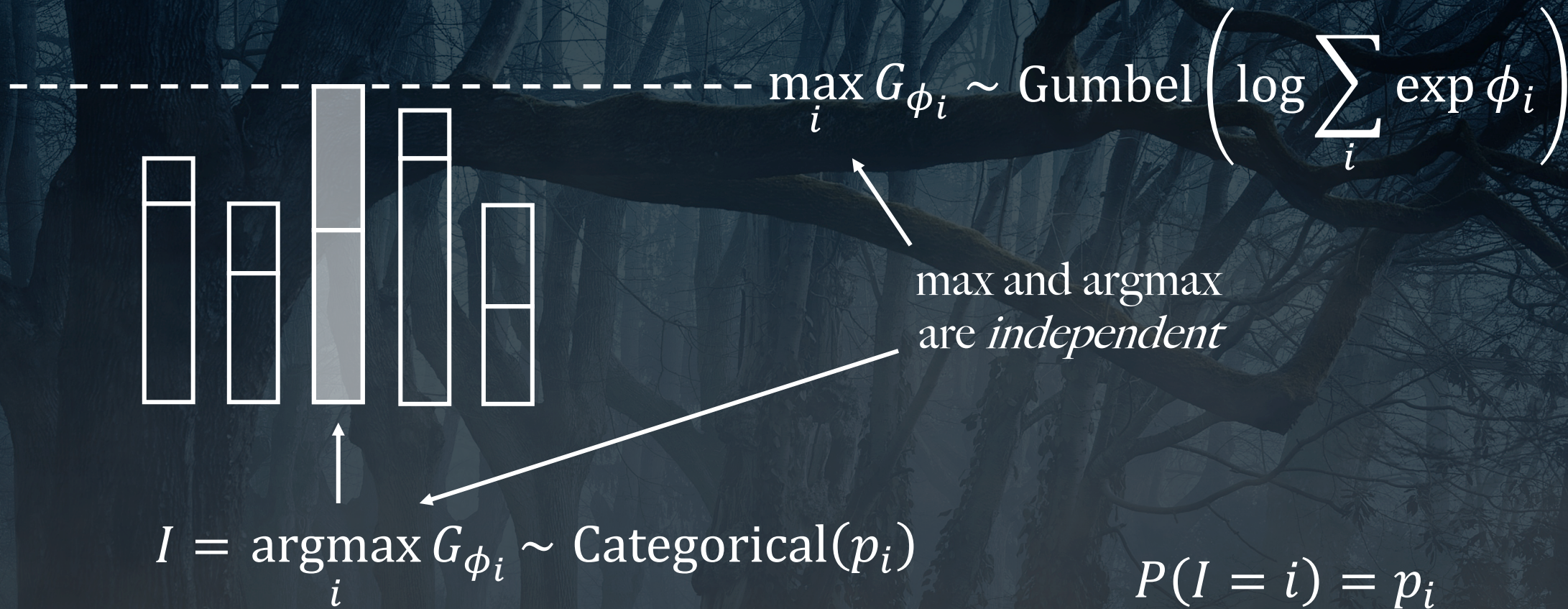
$$G_{\phi_i} \sim \text{Gumbel}(\phi_i)$$

perturbed log-probability

The Gumbel-Max Trick

"Prof. Gumbeldore"

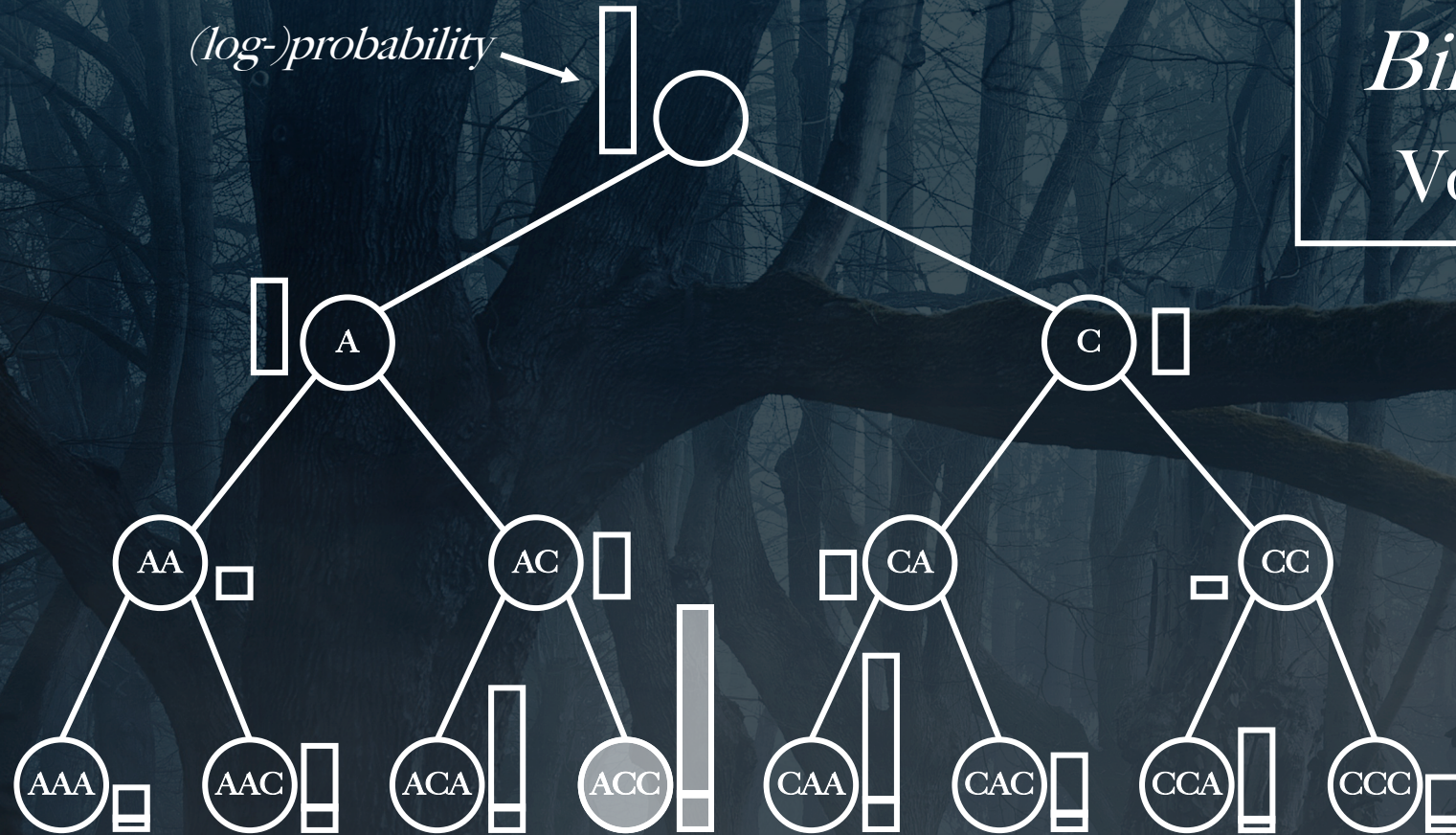
(Gumbel, 1945;
Maddison et al., 2014)



Example

Binarese language model

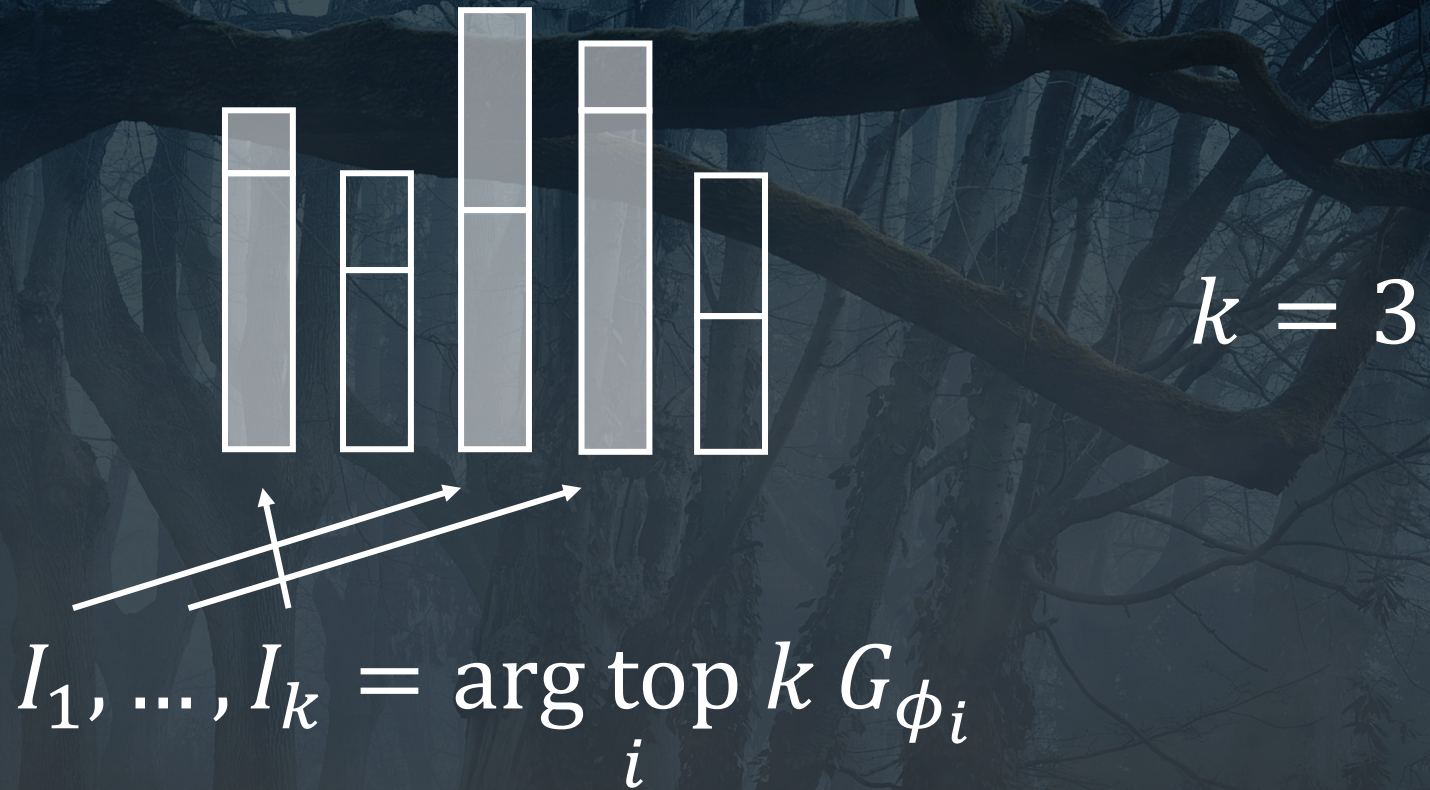
Vocabulary: {**A**_{bra}, **C**_{adabra}}



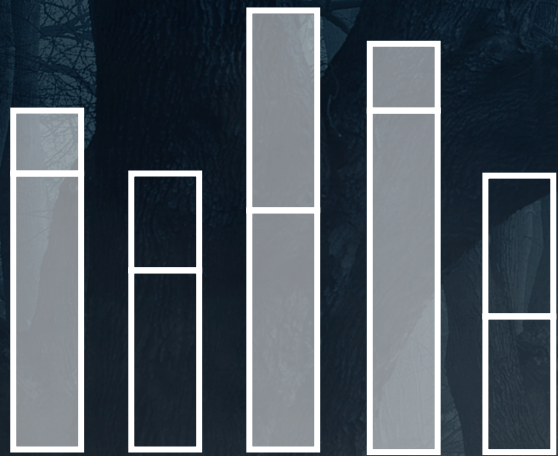
This will be
our sample!

*What if we want
a sample from
our model?*

*What happens if, instead of 1 (one),
we take the k largest elements (top k)?*



The 'Gumbel-Top- k ' Trick



$$I_1, \dots, I_k = \arg \operatorname{top}_k G_{\phi_i}$$

$$\begin{aligned} P(I_1 = i_1, \dots, I_k = i_k) &= p_{i_1} \cdot \frac{p_{i_2}}{1-p_{i_1}} \cdot \dots \cdot \frac{p_{i_k}}{1-\sum_{\ell=1}^{k-1} p_{i_\ell}} \\ &= \prod_{j=1}^k \frac{p_{i_j}}{1-\sum_{\ell=1}^{j-1} p_{i_\ell}} \end{aligned}$$

Also known as
Plackett-Luce

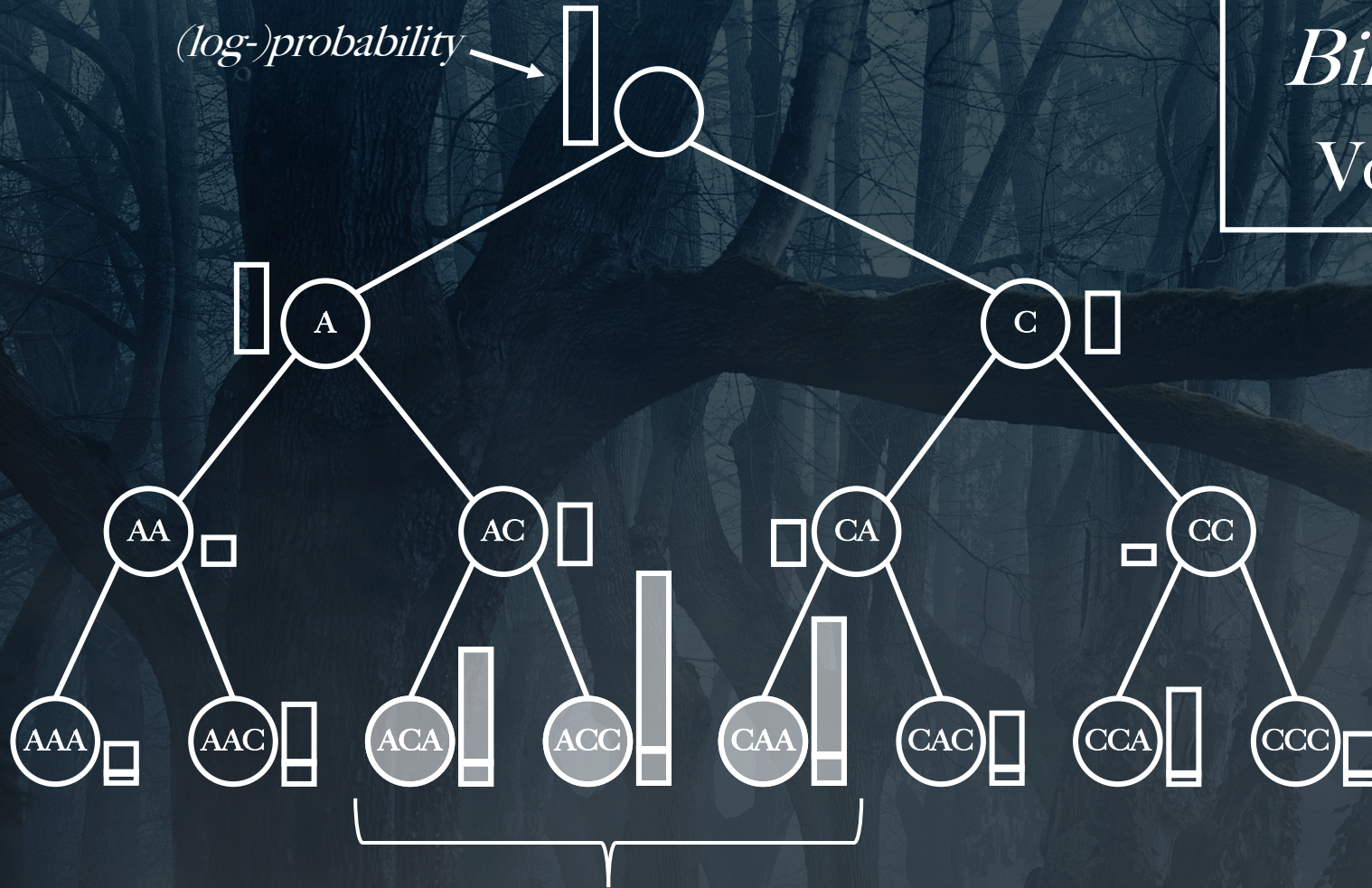
*This is equivalent to repeated
sampling without replacement!*

(Vieira, 2014)

Example

Binarese language model

Vocabulary: $\{A_{bra}, C_{adabra}\}$



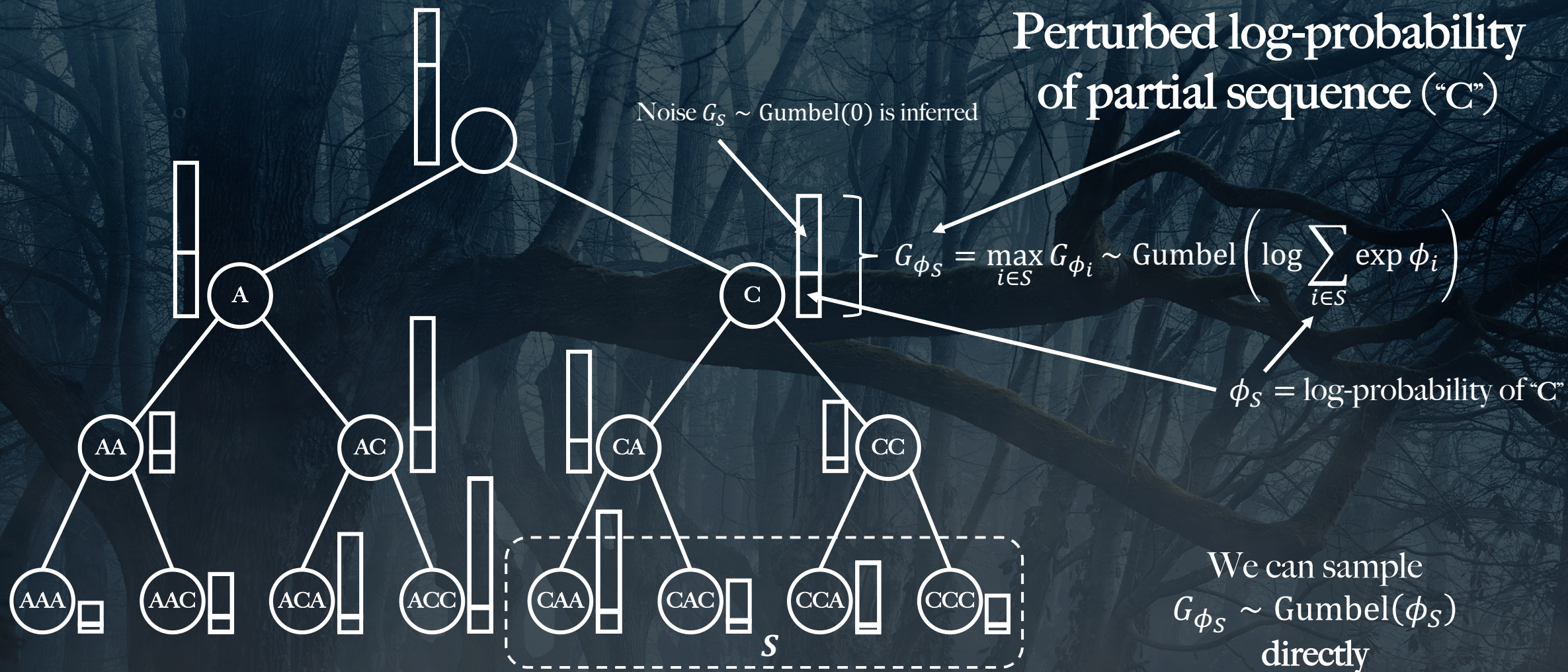
This will be our set of samples!

We can get a set of unique samples from our model!

PROBLEM

In general, constructing
the full tree is not
possible...

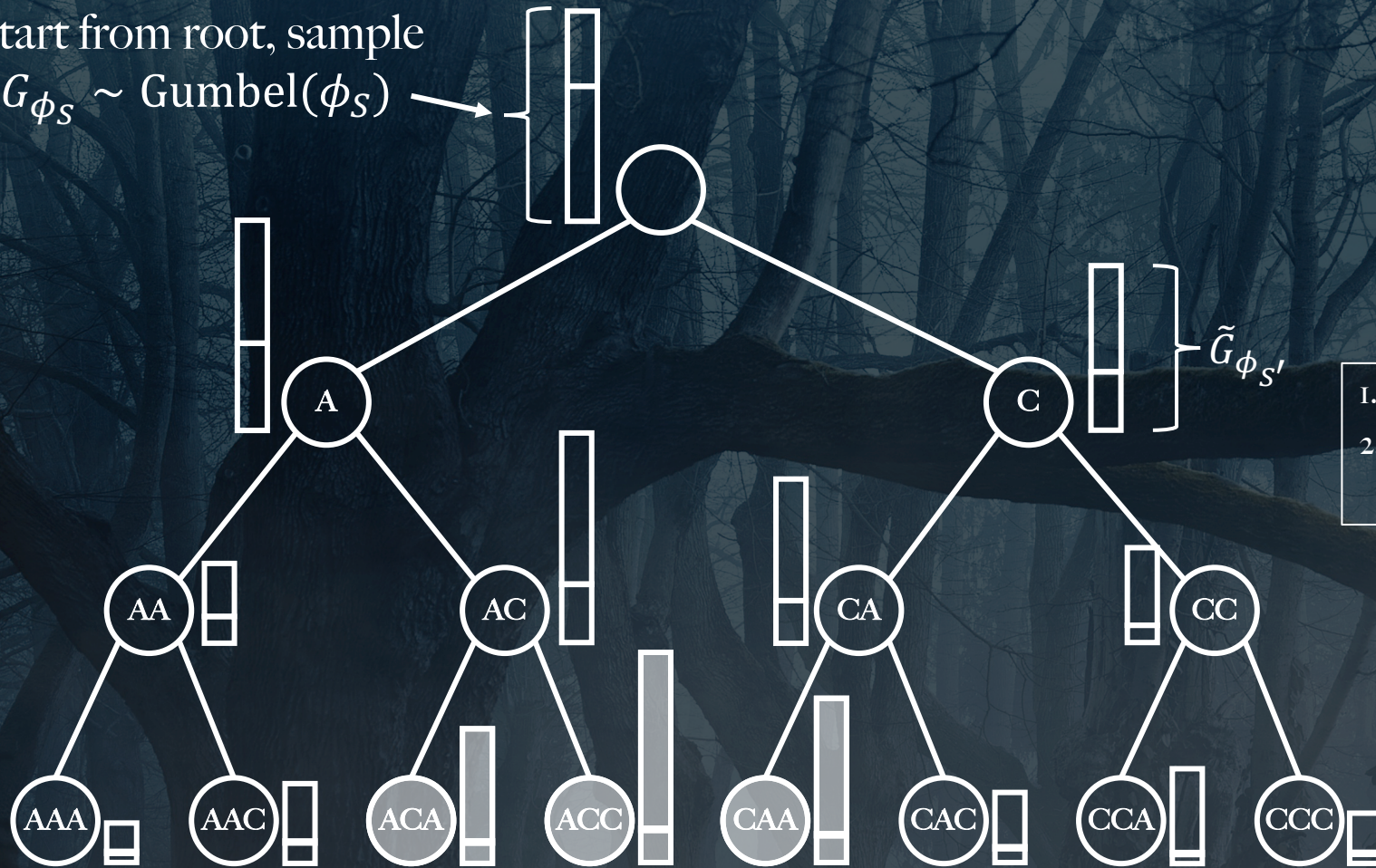
... but we don't have to!



Look at maximum of perturbed log-probabilities in subtree

Start from root, sample

$$G_{\phi_S} \sim \text{Gumbel}(\phi_S)$$



Sample children

$G_{\phi_{S'}}$ conditionally on

$$\max_{S' \in \text{Children}(S)} G_{\phi_{S'}} = G_{\phi_S}$$

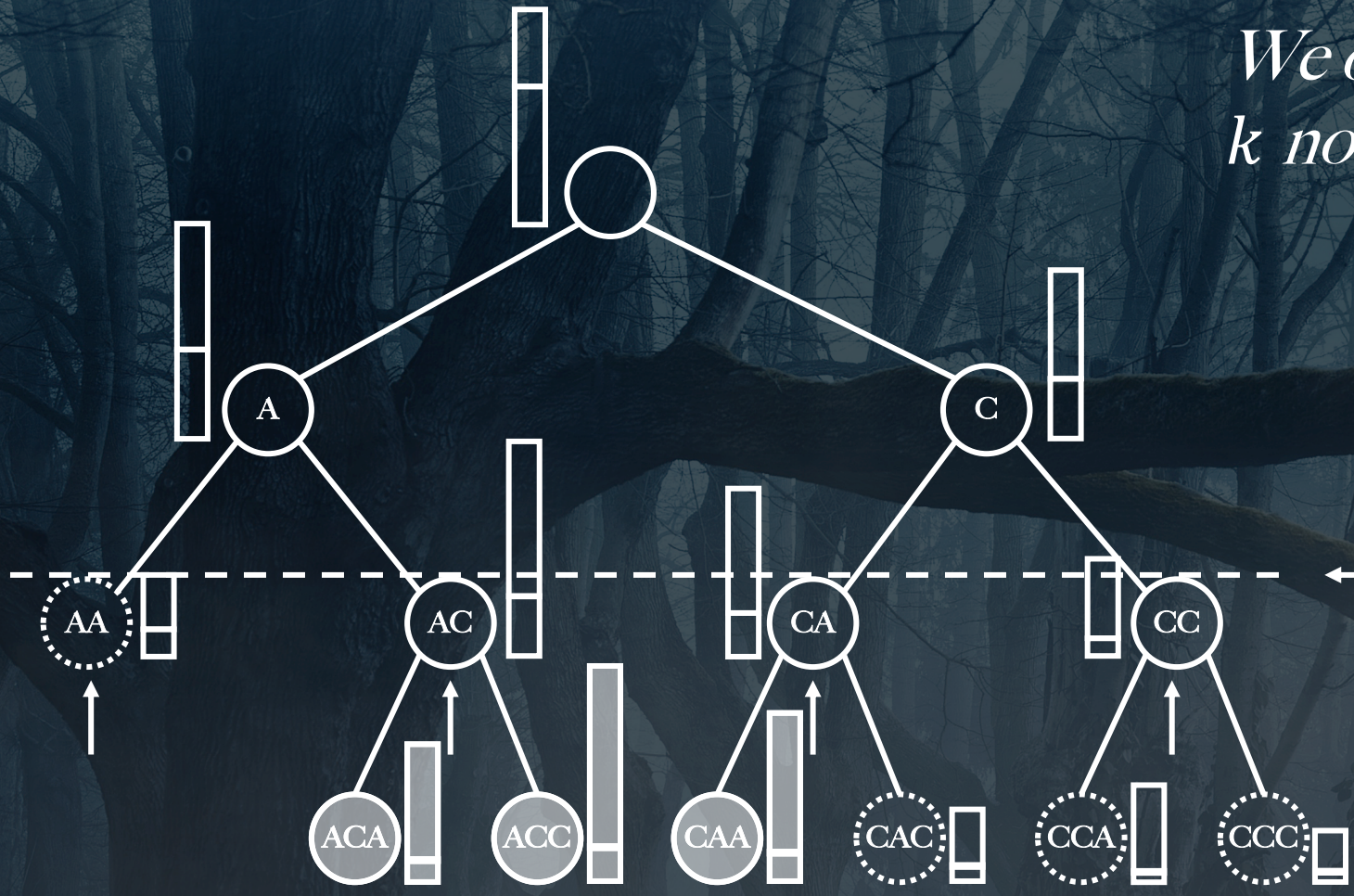
1. sample $G_{\phi_{S'}}$ independently, compute $Z = \max_{S'} G_{\phi_{S'}}$
2. 'shift' Gumbels in (negative) exponential space:
$$\tilde{G}_{\phi_{S'}} = -\log(\exp(-G_{\phi_S}) - \exp(-Z) + \exp(-G_{\phi_{S'}}))$$

... the result is equivalent to sampling G_{ϕ_i} for leaves directly!

Top-down sampling

(Maddison et al., 2014)

We only need to expand the top k nodes at each level in the tree



Threshold

Each top k node generates (at least) one leaf (maximum) above threshold

At least k leaves will be above threshold

Other nodes only generate leaves below threshold

No need to expand



The Key Insight

We only need to expand the top k nodes at each level in the tree



↑
This is a
beam search

Top k according to
perturbed log-probability
= ← Gumbel-Top- k
trick

Sampling (without
replacement)



A white scroll graphic with the text "Stochastic Beam Search" written in a serif font across its center. The scroll is unrolled and has a slight curve.

Stochastic Beam Search

- A beam search that *samples* the nodes to expand
- But... samples children *conditionally* on parent
- The result is a sample without replacement from the full sequence model
- Is a generalization of ancestral sampling ($k = 1$)

Important!







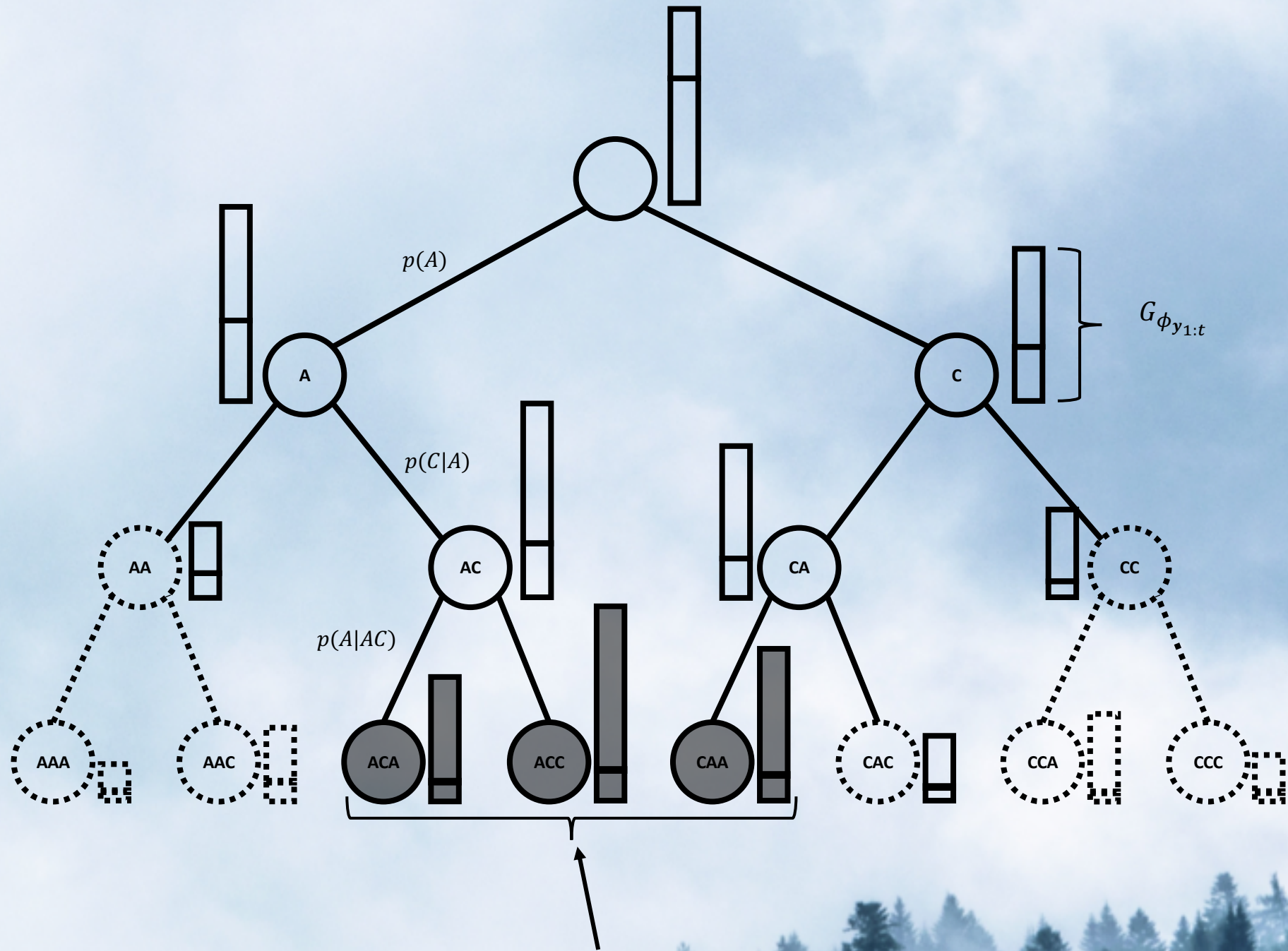
Ancestral Gumbel-Top-*k* Sampling

Ancestral Gumbel-Top- k Sampling

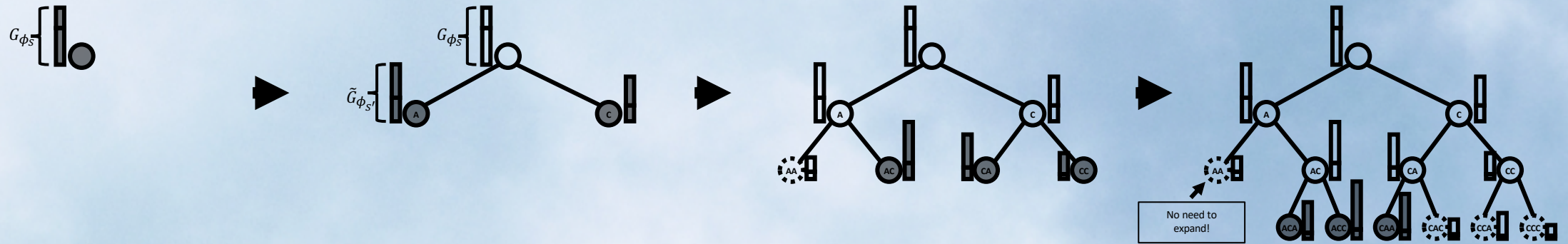
Generalizes Stochastic Beam Search

Expands $1 \leq m \leq k$ nodes per iteration

Applies to discrete valued Bayes networks

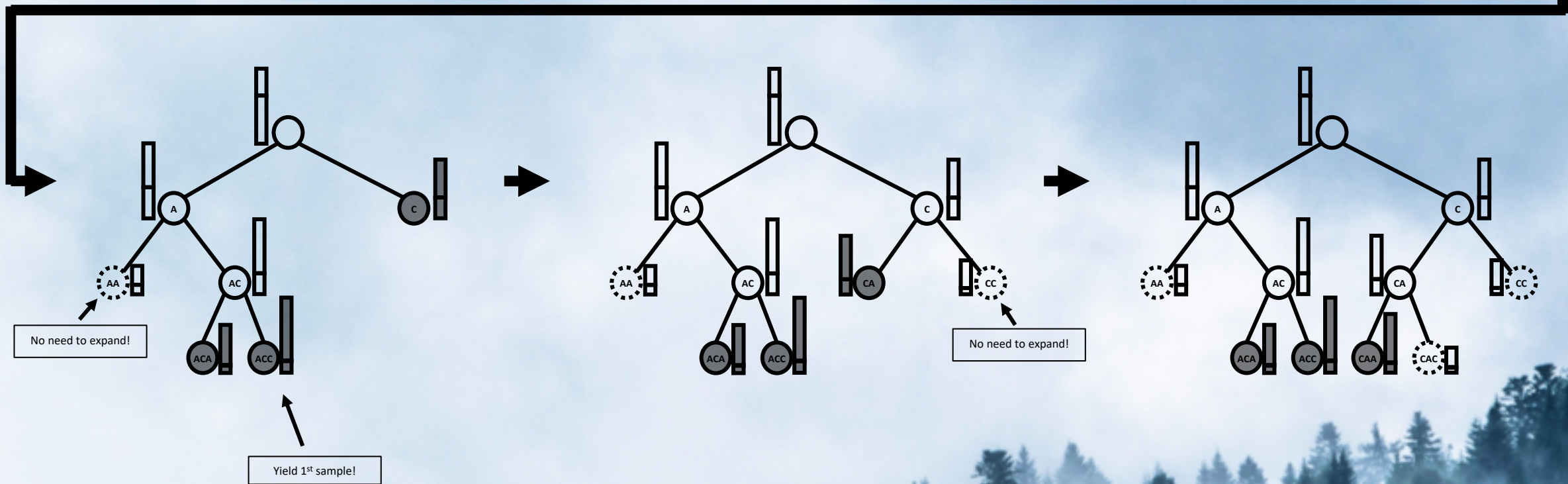
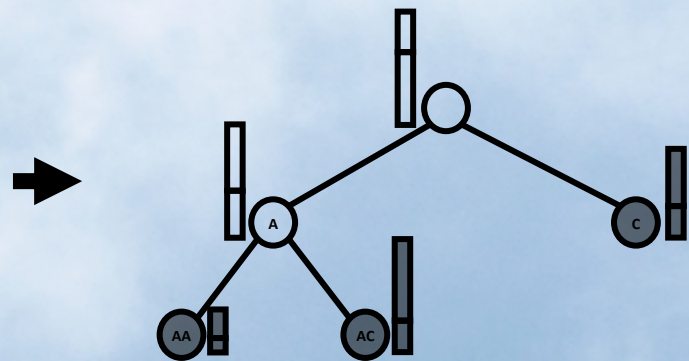
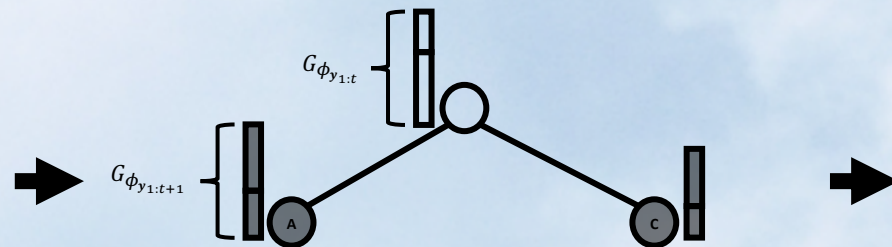
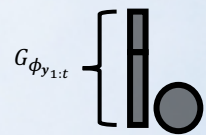


Stochastic Beam Search



$$m = k (= 3)$$

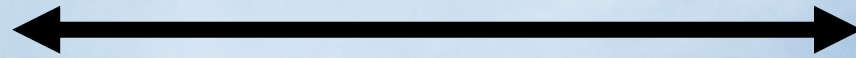
$m = 1$



Ancestral Gumbel-Top- k Sampling

$$m = 1$$

$$m = k$$



Sequential

Parallel

Incremental

Batch

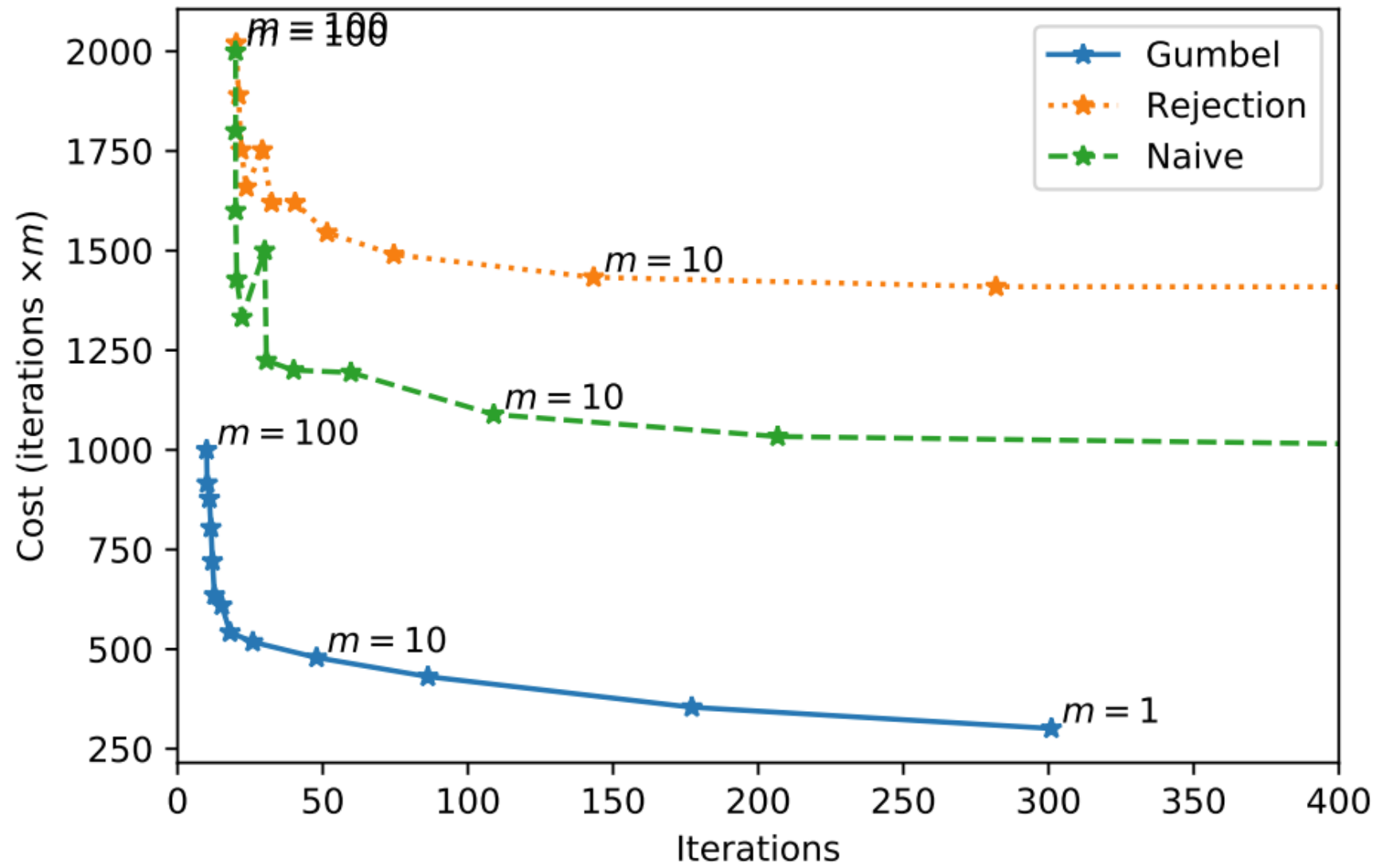
More iterations

Fewer iterations

Less computation

More computation

Cost vs. iterations ($c = 0.5, k = 100$)



Ancestral Gumbel-Top- k Sampling

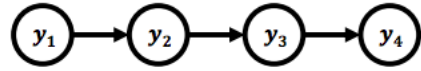
Generalizes Stochastic Beam Search

Expands $1 \leq m \leq k$ nodes per iteration

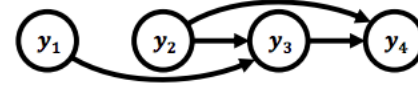
Applies to discrete valued Bayes networks



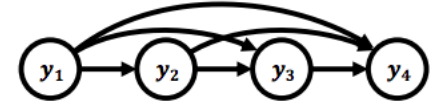
(a) Independent
 $p(\mathbf{y}) = \prod_v p(y_v)$



(b) Markov chain
 $p(\mathbf{y}) = p(y_1) \prod_{t>1} p(y_t | y_{t-1})$

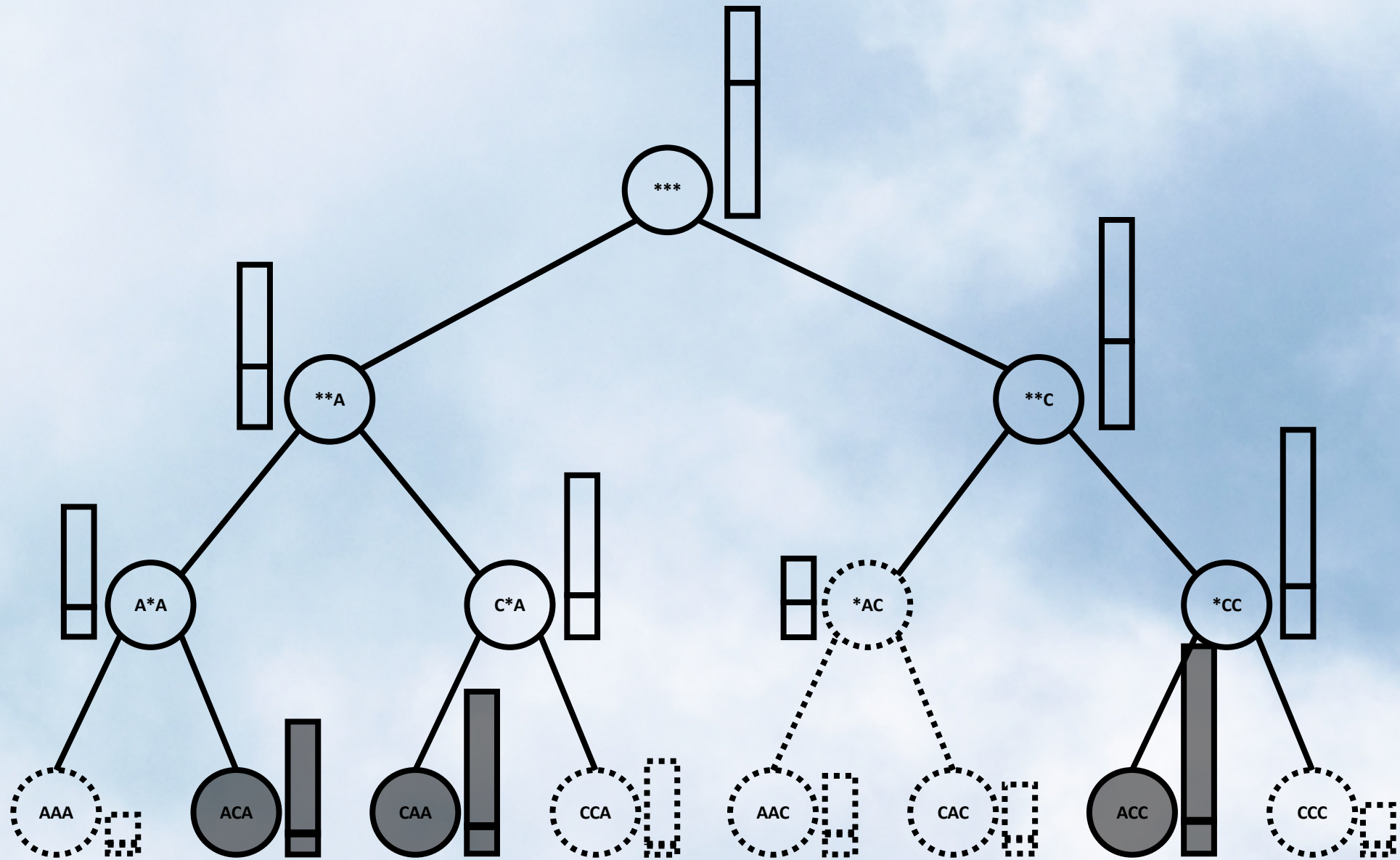


(c) General network
 $p(\mathbf{y}) = \prod_v p(y_v | \mathbf{y}_{\text{pa}(v)})$

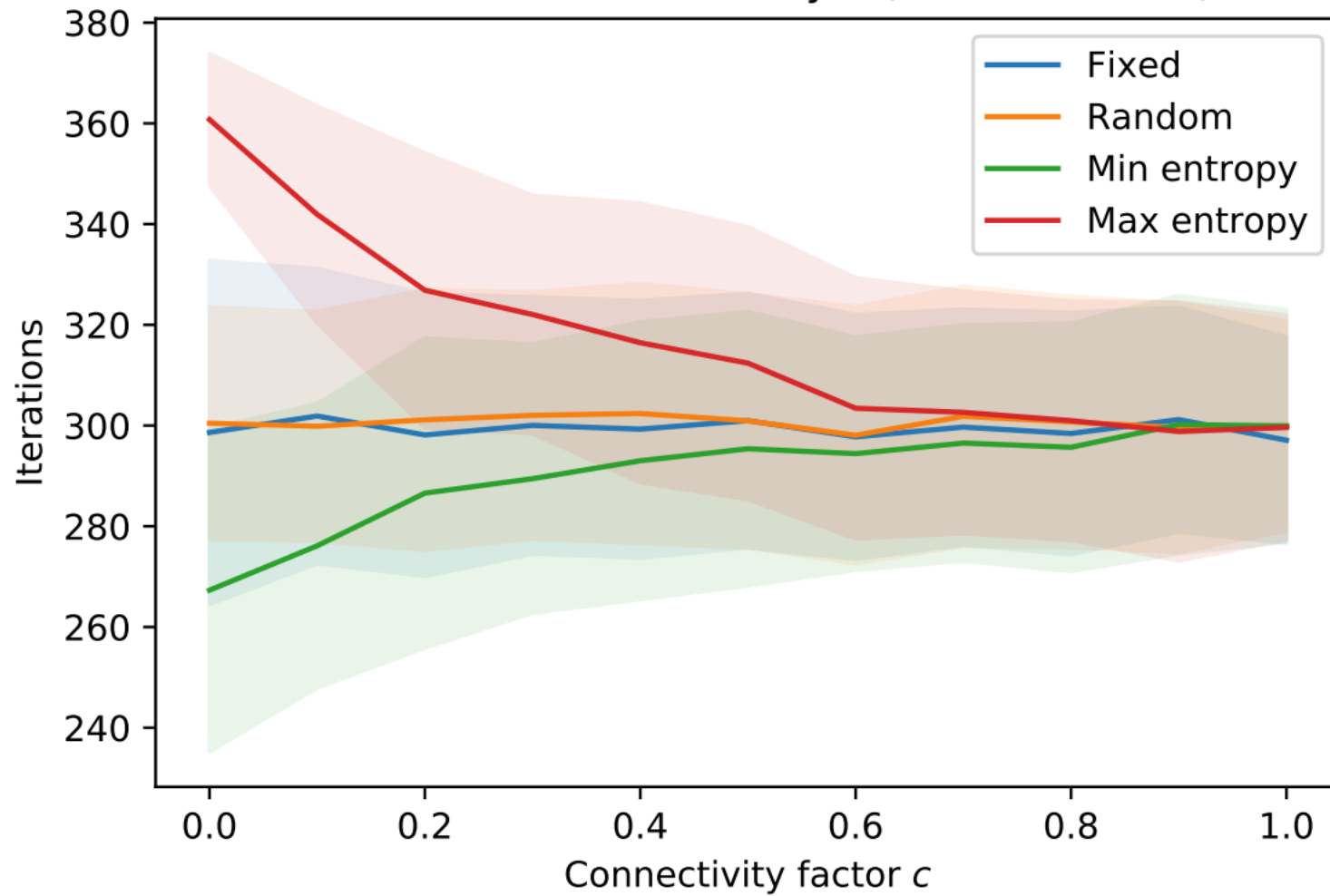


(d) Sequence model
 $p(\mathbf{y}) = \prod_t p(y_t | \mathbf{y}_{1:t-1})$

Figure 1: Examples of Bayesian networks.



Iterations vs. connectivity c ($m = 1, k = 100$)



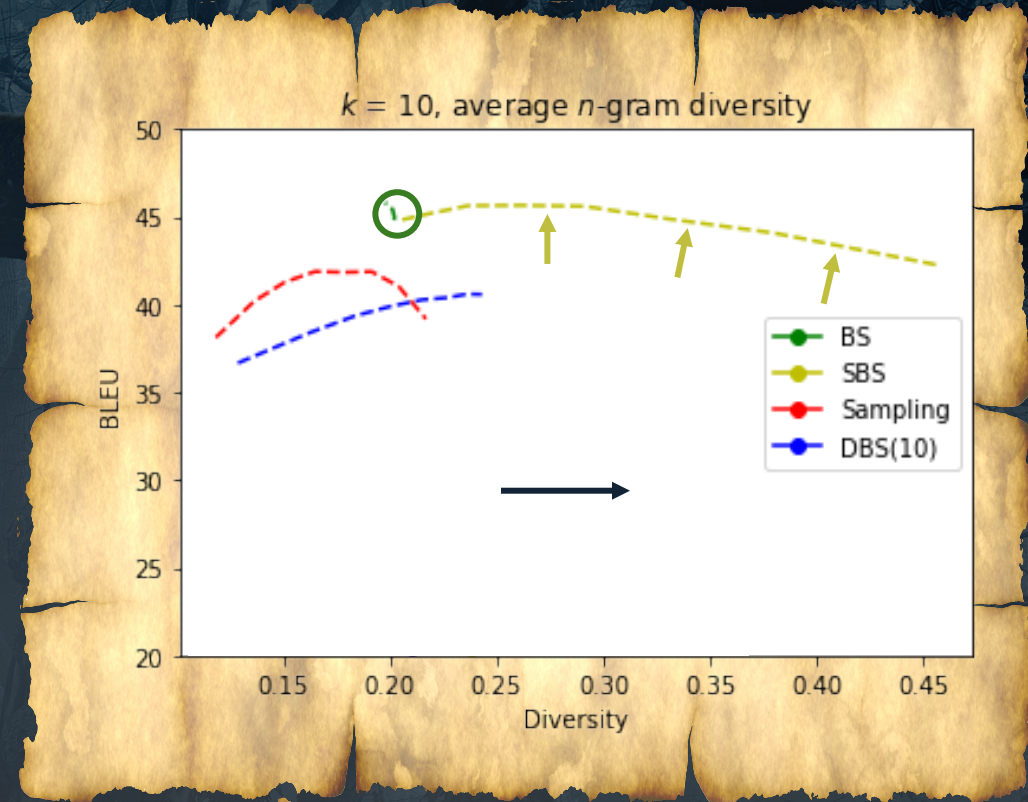
Is this useful?



Experiments

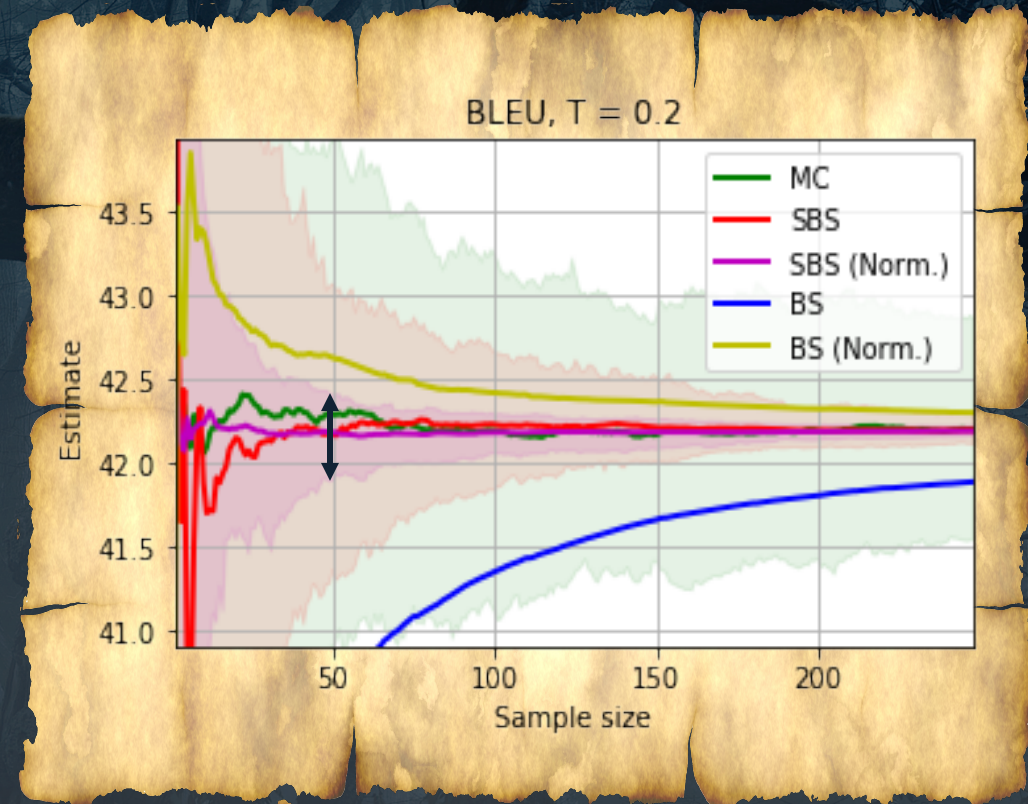
Translation Diversity

- Generate k translations
- Plot BLEU against diversity
- Vary softmax temperature
- Compare:
 - Beam Search
 - Stochastic Beam Search
 - Sampling
 - Diverse Beam Search (Vijayakumar et al., 2018)



BLEU Score Estimation

- Estimate expected sentence-level BLEU
- Plot mean and 95% interval vs. num samples
- Compare:
 - Monte Carlo Sampling
 - Stochastic Beam Search with (normalized) Importance Weighted estimator
 - Beam Search with deterministic estimate





Can we use it
for training?

Estimating Gradients for Discrete Distributions by Sampling Without Replacement

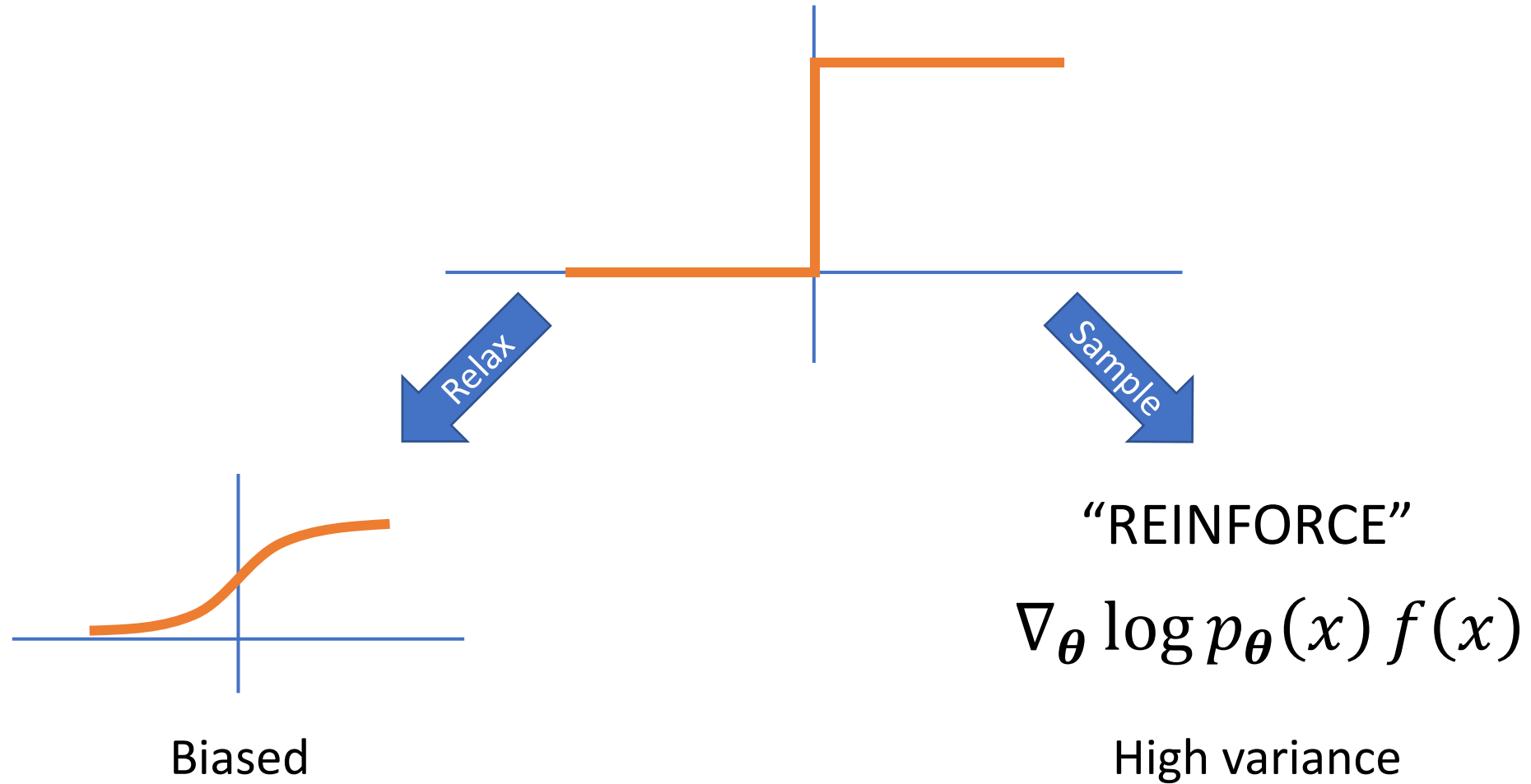
Wouter Kool, Herke van Hoof & Max Welling

International Conference on Learning Representations (ICLR) 2020

Problems of discrete nature

- Reinforcement Learning
- Machine Translation / Image Captioning
- Discrete Latent Variable Modelling
- (Hard) Attention

Gradient of discrete operation



REINFORCE

$$\nabla_{\theta} E_{p_{\theta}(x)}[f(x)] = E_{p_{\theta}(x)}[\nabla_{\theta} \log p_{\theta}(x) f(x)]$$

REINFORCE

$$\nabla_{\theta} E_{p_{\theta}(x)}[f(x)] \approx \nabla_{\theta} \log p_{\theta}(x) f(x)$$

REINFORCE with multiple samples

$$\nabla_{\theta} E_{p_{\theta}(x)}[f(x)] \approx \frac{1}{k} \sum_{i=1}^k \nabla_{\theta} \log p_{\theta}(x_i) f(x_i)$$

REINFORCE with baseline

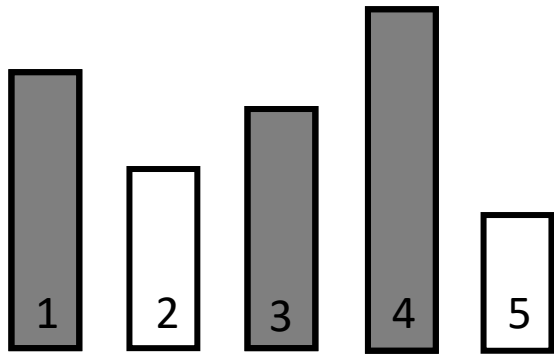
$$\nabla_{\theta} E_{p_{\theta}(x)}[f(x)] \approx \frac{1}{k} \sum_{i=1}^k \nabla_{\theta} \log p_{\theta}(x_i) \left(f(x_i) - \underbrace{\frac{\sum_{j \neq i} f(x_j)}{k-1}}_{\text{Baseline}} \right)$$

Sampling
without
replacement

Since duplicate samples
are uninformative!

*In a deterministic setting

Sampling without replacement



$B = (3, 4, 1)$

$$\begin{aligned} p(B) &= p(b_1) \\ &\times \frac{p(b_2)}{1 - p(b_1)} \\ &\times \frac{p(b_3)}{1 - p(b_1) - p(b_2)} \end{aligned}$$

Ordered samples without replacement

$$p(B) = \prod_{i=1}^k \frac{p(b_i)}{1 - \sum_{j < i} p(b_j)}$$

Sequence $B = (3,4,1)$

Unordered samples without replacement

$$p(B) = \prod_{i=1}^k \frac{p(b_i)}{1 - \sum_{j < i} p(b_j)}$$

Set $S = \{1,3,4\}$

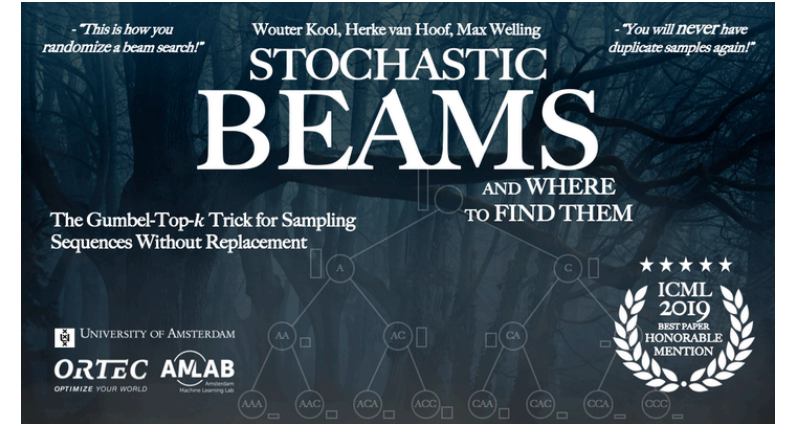
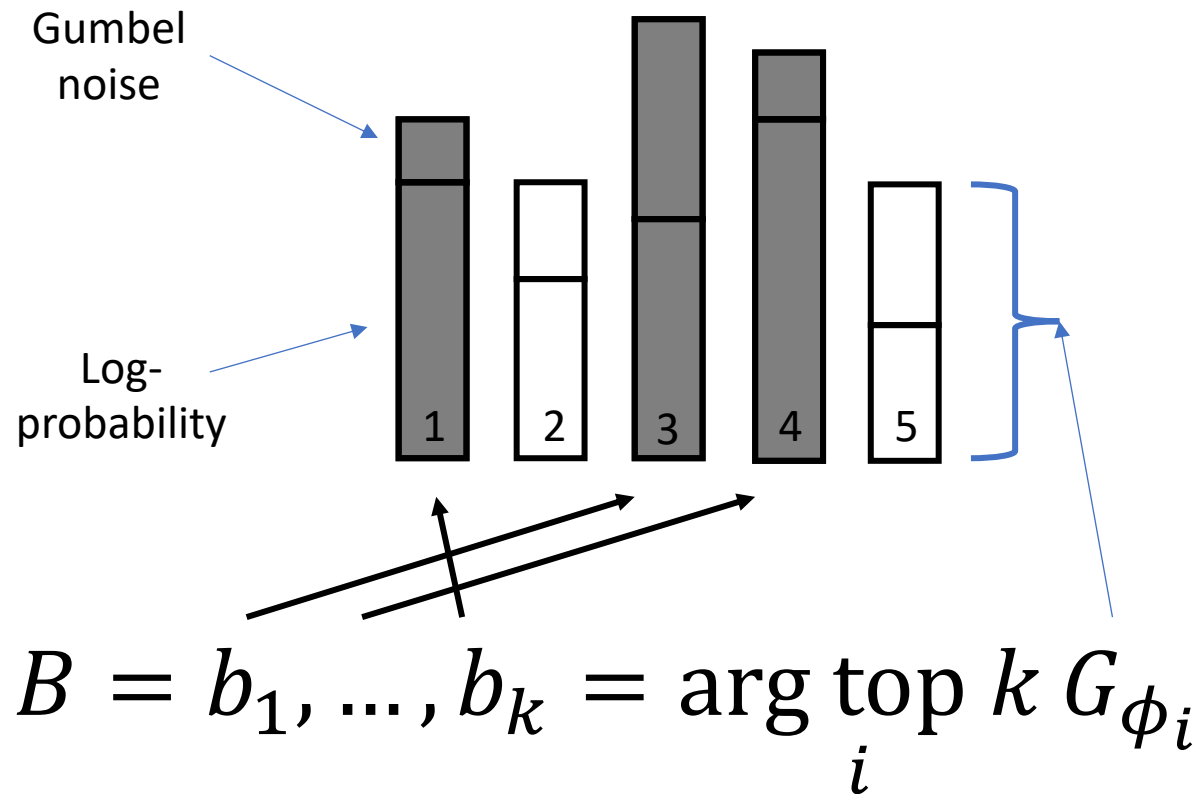
Unordered samples without replacement

$$p(S) = \sum_{B \in \mathcal{B}(S)} p(B) = \sum_{B \in \mathcal{B}(S)} \prod_{i=1}^k \frac{p(b_i)}{1 - \sum_{j < i} p(b_j)}$$

Set $S = \{1,3,4\}$

Sum over $k!$
permutations

Gumbel-Top- k sampling



<https://arxiv.org/abs/1903.06059>

<http://www.jmlr.org/papers/v21/19-985.html>

$$B = (3, 4, 1)$$
$$S = \{1, 3, 4\}$$

Back to our problem

$$\nabla_{\boldsymbol{\theta}} E_{p_{\boldsymbol{\theta}}(x)} [f(x)]$$

Estimating the expectation

$$E_{p_{\theta}(x)}[f(x)]$$

The single sample estimator

$$E_{p_{\theta}(x)}[f(x)] = E_{p_{\theta}(B)}[f(b_1)]$$

Separating the expectation

$$E_{p_{\theta}(x)}[f(x)] = E_{p_{\theta}(S)} \left[E_{p_{\theta}(B|S)}[f(b_1)] \right]$$

↑
Set of
unordered
samples

↑
Conditional
distribution of
their order

Separating the expectation

$$E_{p_{\theta}(x)}[f(x)] = E_{p_{\theta}(s)} \left[E_{p_{\theta}(b_1|s)}[f(b_1)] \right]$$

Rao-Blackwellizing the estimator

$$E_{p_{\theta}(x)}[f(x)] = E_{p_{\theta}(S)} \left[E_{p_{\theta}(b_1|S)}[f(b_1)] \right]$$

$$E_{p_{\theta}(b_1|S)}[f(b_1)] = \sum_{s \in S} P(b_1 = s | S) f(s)$$

Rao-Blackwellizing the estimator

$$E_{p_{\theta}(x)}[f(x)] = E_{p_{\theta}(S)} \left[E_{p_{\theta}(b_1|S)}[f(b_1)] \right]$$

$$E_{p_{\theta}(b_1|S)}[f(b_1)] = \sum_{s \in S} P(b_1 = s|S) f(s)$$

$$P(b_1 = s|S) = \frac{P(S|b_1 = s)P(b_1 = s)}{P(S)}$$

Rao-Blackwellizing the estimator

$$E_{p_{\theta}(x)}[f(x)] = E_{p_{\theta}(S)} \left[E_{p_{\theta}(b_1|S)}[f(b_1)] \right]$$

$$E_{p_{\theta}(b_1|S)}[f(b_1)] = \sum_{s \in S} P(b_1 = s|S) f(s)$$

$$P(b_1 = s|S) = \underbrace{\frac{P(S|b_1 = s)}{P(S)}}_{\text{Leave-one-out ratio } R(S, s)} \underbrace{P(b_1 = s)}_{p_{\theta}(s)}$$

Leave-one-out ratio $R(S, s)$

$p_{\theta}(s)$

Rao-Blackwellizing the estimator

$$E_{p_{\theta}(x)}[f(x)] = E_{p_{\theta}(S)} \left[E_{p_{\theta}(b_1|S)}[f(b_1)] \right]$$

$$E_{p_{\theta}(b_1|S)}[f(b_1)] = \sum_{s \in S} P(b_1 = s|S) f(s)$$

$$P(b_1 = s|S) = R(S, s)p_{\theta}(s)$$

Rao-Blackwellizing the estimator

$$E_{p_{\theta}(x)}[f(x)] = E_{p_{\theta}(S)} \left[E_{p_{\theta}(b_1|S)}[f(b_1)] \right]$$

$$E_{p_{\theta}(b_1|S)}[f(b_1)] = \sum_{s \in S} R(S, s) p_{\theta}(s) f(s)$$

Rao-Blackwellizing the estimator

$$E_{p_{\theta}(x)}[f(x)] = E_{p_{\theta}(S)} \left[\underbrace{\sum_{s \in S} R(S, s) p_{\theta}(s) f(s)}_{\text{Unordered set estimator}} \right]$$

Combining with REINFORCE

$$E_{p_{\theta}(x)} [f(x)]$$
$$= E_{p_{\theta}(S)} \left[\sum_{s \in S} R(S, s) p_{\theta}(s) f(s) \right]$$

Combining with REINFORCE

$$E_{p_{\theta}(x)} [\nabla_{\theta} \log p_{\theta}(s) f(x)]$$
$$= E_{p_{\theta}(S)} \left[\sum_{s \in S} R(S, s) p_{\theta}(s) \nabla_{\theta} \log p_{\theta}(s) f(s) \right]$$

Combining with REINFORCE

$$\begin{aligned}\nabla_{\theta} E_{p_{\theta}(x)} [f(x)] &= E_{p_{\theta}(x)} [\nabla_{\theta} \log p_{\theta}(s) f(x)] \\ &= E_{p_{\theta}(s)} \left[\sum_{s \in S} R(S, s) \underbrace{p_{\theta}(s) \nabla_{\theta} \log p_{\theta}(s)}_{\nabla_{\theta} p_{\theta}(s)} f(s) \right]\end{aligned}$$

Combining with REINFORCE

$$\nabla_{\theta} E_{p_{\theta}(x)} [f(x)]$$

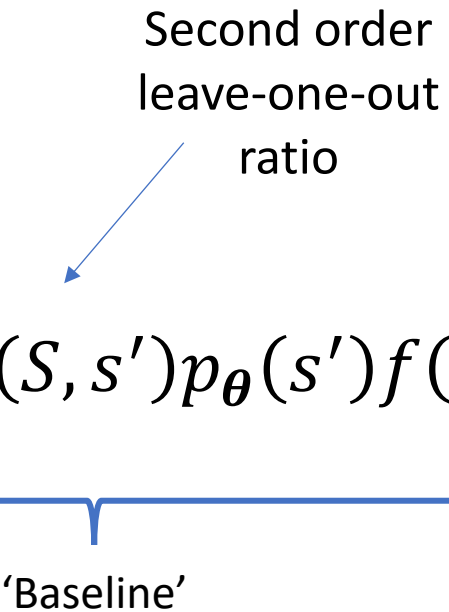
$$= E_{p_{\theta}(S)} \left[\underbrace{\sum_{s \in S} R(S, s) \nabla_{\theta} p_{\theta}(s) f(s)}_{\text{Unordered set policy gradient estimator}} \right]$$

Unordered set policy gradient estimator

Include a baseline

$$\nabla_{\theta} E_{p_{\theta}(x)} [f(x)]$$
$$= E_{p_{\theta}(S)} \left[\sum_{s \in \mathcal{S}} R(S, s) \nabla_{\theta} p_{\theta}(s) \left(f(s) - \underbrace{\sum_{s' \in \mathcal{S}} R^{\setminus s}(S, s') p_{\theta}(s') f(s')}_{\text{'Baseline'}} \right) \right]$$

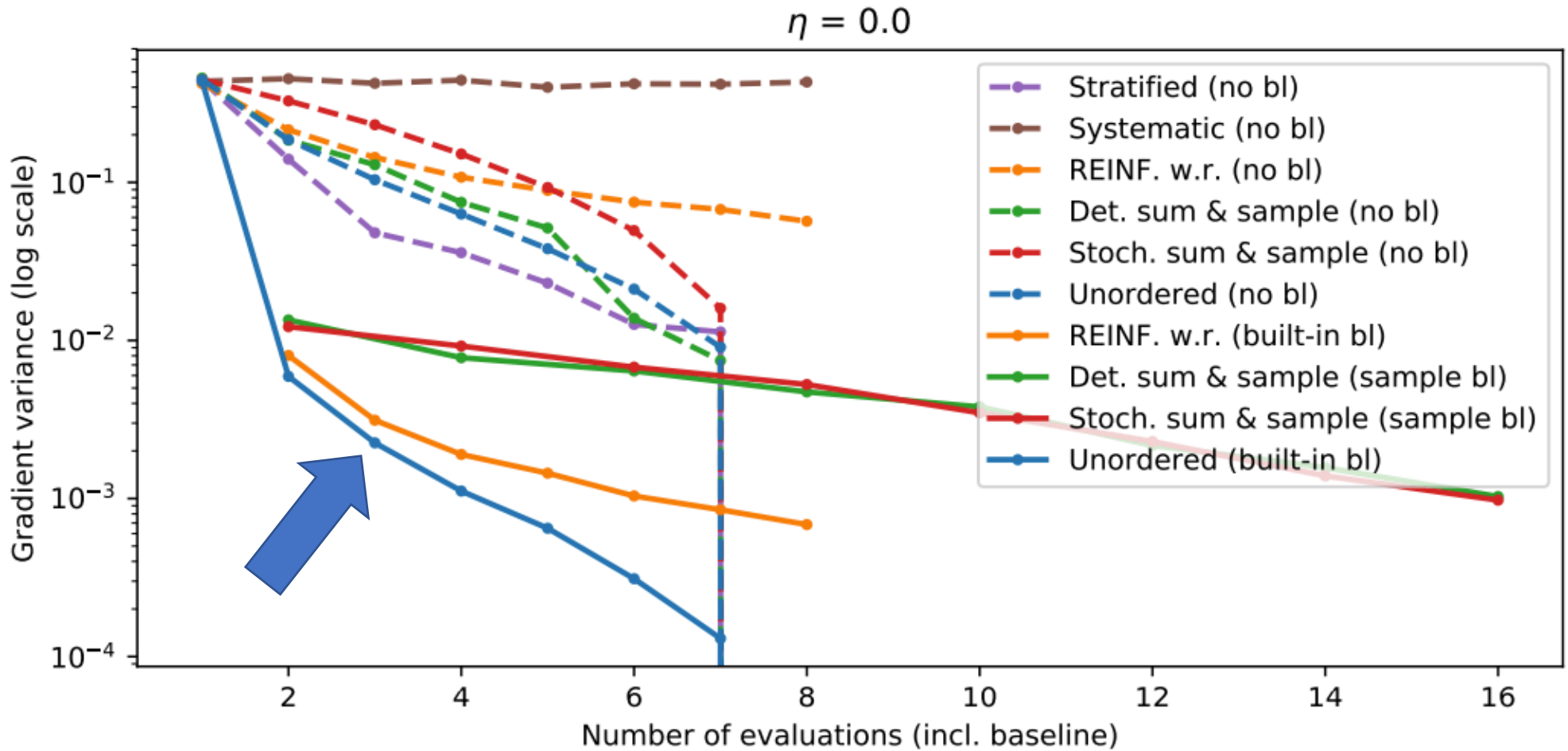
Second order
leave-one-out
ratio



Unbiased!

Experiments

Bernoulli gradient variance



(a) High entropy ($\eta = 0$)

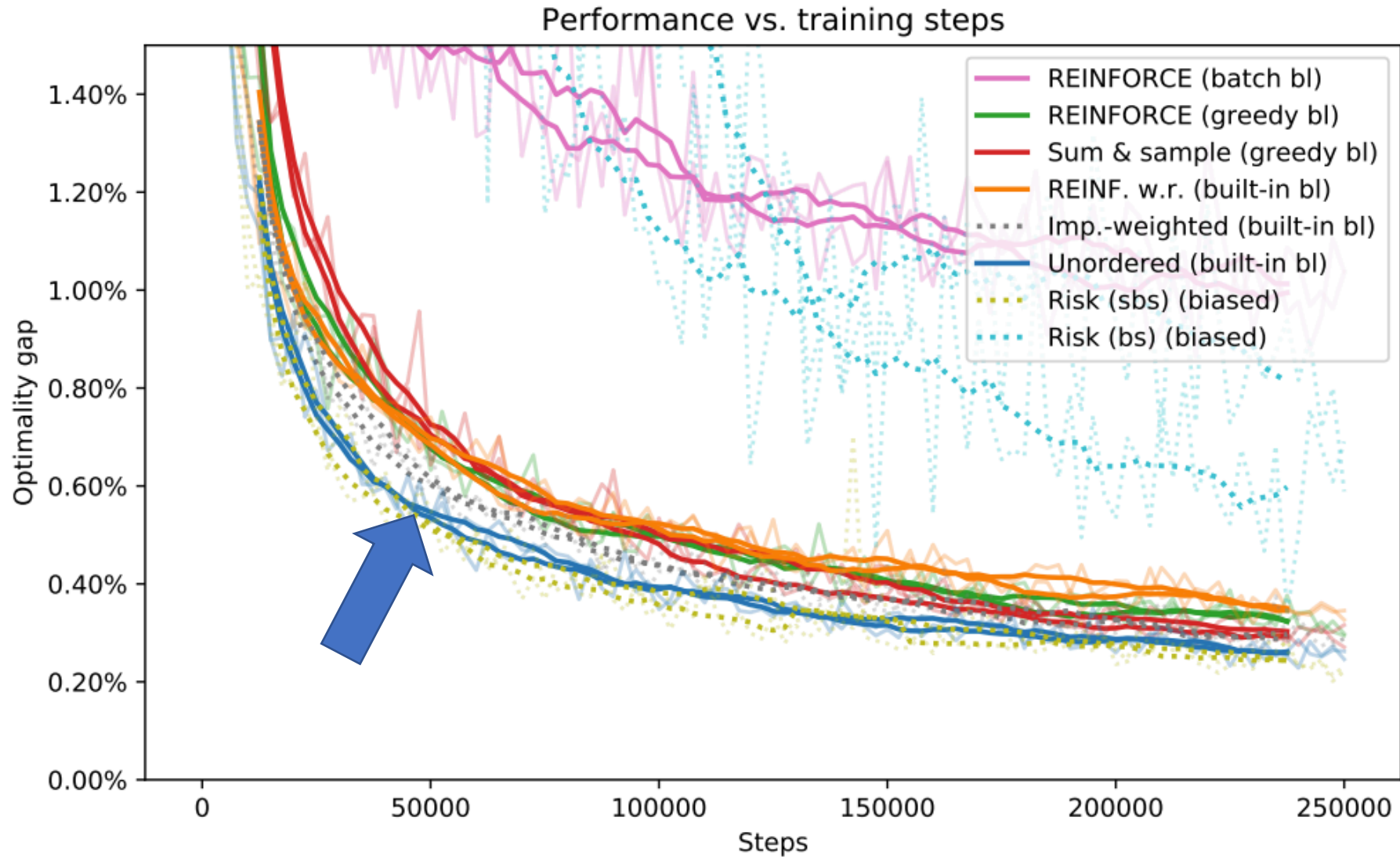
Categorical Variational Auto-Encoder (grad. var.)

Table 1: VAE gradient log-variance of different unbiased estimators with $k = 4$ samples.

Domain	ARSM	RELAX	REINFORCE		Sum & sample		REINF. w.r.	Unordered
			(no bl)	(sample bl)	(no bl)	(sample bl)	(built-in bl)	(built-in bl)
Small 10^2	13.45	11.67	11.52	7.49	6.29	6.29	6.65	6.29
Large 10^{20}	15.55	15.86	13.81	8.48	13.77	8.44	7.06	7.05



Travelling Salesman Problem



Take away

The unordered set estimator

- Low-variance
- Unbiased
- Alternative to Gumbel-Softmax

End of story